

GPU Accelerated Neural Networks With Dynamic Topology

Travis Chung
Project Leader

Shashank Golla
Code Architecture

Suhil Sam Kiswani
Simulation Details

April 24, 2015

Todo list

Abstract

Neural Networks lend themselves naturally to being parallelized, although certain caveats exist. Furthermore, there exists little work on the benefits of parallelizing networks having a dynamic topology. We demonstrate a parallel implementation of the the Synfire-Growth simulation (developed by Jun and Jin [?]) that runs on the GPU. Despite the dynamic topology of the neural network, our results show that porting the work of [?] to the GPU provided useful runtime performance benefits.

1 Introduction

The human brain is one of the most complicated and interesting biological systems in nature. Despite our poor understanding of the brain's complex structure, it is common to describe it as a biological computer, consisting of an dense and tangled network of neurons, which serve as pathways for electrical current. This analogy provides a useful starting point for computational models, and is often employed in the machine learning domain.

Neural Networks (or NN's) can be understood most simply as a complete weighted graph, with vertices representing neurons, and weighted edges representing the strength of a neural connection. This model can be enhanced by allowing the weights to change in response to external stimuli, representing a dynamic topology of connected neurons.

This paper describes the simulation of neural networks with dynamic topology on the GPU using the model developed by [?] in 2007. We describe how the implementation of [?] was refactored to run on the GPU and take advantage of parallelization. Our work concludes with a cost-benefit analysis of running this simulation on the GPU, and how neural networks naturally lend themselves to the parallel environment.

Remove
the todo
list and all
todos.

double
check this,
fill in some
actual
data as re-
sults come
in. add a
tiny bit
more.

elaborate
and make
sure its
100% obvi-
ous our
imple-
mentation
is just a
refactoring
of [?]

2 Background

Primarily, the study the neuroscience is concerned with the behavior of the brain, the way it is structured, the way it learns, how it develops, how it adapts, and changes with respect to stimuli. While understanding of the central nervous system continues to grow, there is so much that is unknown about the brain: the storage and accessing of memories, how the brain retains information, the idea of consciousness, and how sensory input is translated into smell, taste, or pain. The modeling and simulating of brain activity is vital in observing the behavior of the brain and building intuition.

There are generally two families of algorithms for the simulation of neural networks, described in detail by [?]. The two families are synchronous (“clock-driven”) or asynchronous (“event-driven”) algorithms. In synchronous algorithms neurons are updated only when they receive or emit a spike, whereas “clock-driven” algorithms update all neurons simultaneously at every tick of a clock. There are plenty of simulations using synchronous algorithms, because the spike times are a defined time grid. To get exact simulations of neuron spiking, asynchronous algorithms are recommended. Asynchronous algorithms have been developed for simpler models, but for very large networks, the simulation time and number of spikes becomes problematic for computation.

3 Previous Work

Despite the group’s lack of domain knowledge regarding neural networks, the work of [?] provided the initial foundation and direction for our work. The comparisons of the Hodgkin-Huxel (HH) and Izhikevich models of neural networks provided by [?] was a significant aide in deciding which model to use for our simulation. In addition to their analysis, [?] provides evidence for parallel implementations on the GPU that achieve speedups of greater than 110 times their corresponding CPU implementations. Though we could not reproduce a speedup of their magnitude, it gave our work its initial direction.

The model we chose to implement on the GPU was one developed and detailed by [?], suggested by an informal advisor for this project.

In addition to developing the simulation model, it was also implemented by [?]. By using their implementation as our foundation, we were able to run a faithful representation of their work on the GPU which provides a convenient metric for the efficacy of our work.

4 Implementation Details

synfire refactoring description: improvements, complications etc. short and sweet

proposal material. clean up, rephrase statements, and trim the fat.

requires a clean-up and a couple of revisions to fit the tone of the paper. Include a paragraph on the synfire paper.

be more specific about how [?] served as a foundation for the project. as it stands it seems out of place and kinda useless. also phrasing is awkward.

expand on model of [?]. give

Algorithm 4.1 Synfire Growth Trial

```
while  $t < \text{TRIAL\_TIME}$  do
  MEMBRANEPOTENTIAL( $dt$ )
  for  $n \in N$  that spiked do                                     ▷ Spike Loop
    Check to see if spiking neuron is saturated.
    Spiking neuron isn't saturated, send spikes along active connections.
    Synaptic Plasticity.
  end for
  for  $n \in N$  that spiked do                                     ▷ Inhibition Phase
    Calculate inhibition.
  end for
   $t \leftarrow t + dt$ 
end while
SYNAPTICDECAY()
```

4.1 On the CPU

depending on time, mention how we tried to help the CPU being competitive in our analysis

4.2 On the GPU

describe some potential + actual optimizations by name; choose 2/3 and go into explicit detail in the subsubsects.

reasons for choosing aforementioned optimizations goes here

4.2.1 Membrane Potential Layer

Due to the fact that each neuron is integrated over dt individually, the membrane potential layer update step lends itself very naturally to being parallelized.

4.2.2 Optimization 2

describe one method of the GPU implementation, and include some of the timings. maybe compare it to CPU? dont say its the best here though

4.2.3 Optimization 3

might be able to get away with only 2 but this serves as a reminder to shoot for the stars

5 Results

IMPORTANT!!

talk about the pros/cons of certain methods and do the major pro/con chart between methods.

be specific about the differences in implementation details. how optimizations can be combined or something to fill up the space

a conclusory sentence that sums up the report

References

- [1] Jun, Joseph K and Jin, Dezhe Z, *Development of neural circuitry for precise temporal sequences through spontaneous activity, axon remodeling, and synaptic plasticity*. PLoS ONE (2007) 2(8): e723. doi: 10.1371/journal.pone.0000723
- [2] Brette R, Rudolph M, Carnevale T, et al. *Simulation of networks of spiking neurons: A review of tools and strategies*. *Journal of computational neuroscience*. 2007;23(3):349-398. doi:10.1007/s10827-007-0038-6.
- [3] Fidjeland, A.K., Shanahan, M.P. *Accelerated simulation of spiking neural networks using GPUs*. The 2010 International Joint Conference on Neural Networks (IJCNN). (2010)